

Databases & Data Mining

Erwin M. Bakker & Stefan Manegold

e.m.bakker

s.manegold

@liacs.leidenuniv.nl

<https://homepages.cwi.nl/~manegold/DBDM/>

<http://liacs.leidenuniv.nl/~bakkerem2/dbdm/>



DBDM: Overview

Period: September 11th - December 4th 2018 (Tuesdays)

Place: Room 312 (LIACS, Snellius building, Niels Bohrweg 1, 2333 CA Leiden)

Time: 15.30 - 17.15

ECTS: 6

Description:

The course Databases & Data Mining consists of a series of lectures in which advanced database and data mining techniques will be discussed, with applications to bioinformatics.

Grading:

There will be 2 database and 2 data mining assignments, i.e., 4 assignments in total, and a final exam (open book). The final grade will be based on a weighted average of the grades obtained for assignments P1, P2, P3, P4 and the Exam (E >5):

Final Grade = $(0.5 * P1 + P2 + 0.5 * P3 + P4 + 3 * E) / 6$.

DBDM: “Registration”

Please send an email

To: s.manegold@liacs.leidenuniv.nl

Subject: [DBDM-2018] Registration

containing the following information:

- Your full name
- Your email address
- Your student ID
- Your affiliation (university)
- Your program / subject

By Sunday 16 September 2018, 23:59 CEST.

<https://homepages.cwi.nl/~manegold/DBDM/>

<http://liacs.leidenuniv.nl/~bakkerem2/dbdm/>

DBDM: (tentative) Schedule

Date	Room	Subject (tentative)	Topic & Lecturer
11-09	312	Introduction	Databases and Data Management for Data Mining Stefan Manegold
18-09	312	Database Technology	
25-09	312	Database Technology	
02-10	312	Data Preprocessing	
09-10	312	No class	
16-10	312	Data Warehousing and OLAP	
23-10	312	Data Cube Technology	Data Mining Techniques and Applications Erwin Bakker
30-10	312	Basic Data Mining Algorithms I	
06-11	312	Basic Data Mining Algorithms II	
13-11	312	Advanced Data Mining Algorithms	
20-11	312	Mining in Bio-Data	
27-11	312	Graph Mining I	
04-12	312	Graph Mining II	

DBDM: Assignments

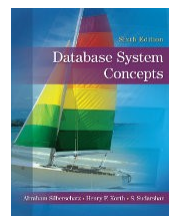
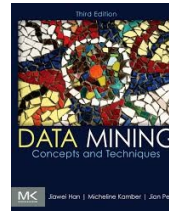
- 2 database assignments & 2 data mining assignments
- Will be announced individually during lectures and posted on website

<https://homepages.cwi.nl/~manegold/DBDM/>

<http://liacs.leidenuniv.nl/~bakkerem2/dbdm/>

DBDM: Recommended Books

- **Data Mining:**
 - J. Han, M. Kamber, J. Pei. **Data Mining Concepts and Techniques (3rd Edition)**, Morgan Kaufman Publishers, July 2011 (ISBN 978-0123814791)
- **Database systems (e.g.):**
 - Ramakrishnan, Gehrke: **Database Management Systems (3rd International Edition)**, McGraw-Hill, 2003 (ISBN 0-07-246563-8)
 - A. Silberschatz, H. F. Korth, S. Sudarshan: **Database System Concepts (6th Edition)**, McGraw-Hill, 2010 (ISBN 0-07-352332-1)



<https://homepages.cwi.nl/~manegold/DBDM/>

<http://liacs.leidenuniv.nl/~bakkerem2/dbdm/>

DBDM: Exam

- **open book exam:** you can take with you your book, and printed course notes (slides). *No electronic equipment is allowed, though.*
- **Materials to be studied:**
 - All content covered and discussed during lectures (slides will be shared).
 - More to be announced.
- **Date:** Monday, January 7, 2019
- **Time:** 14:00 - 17:00
- **Place:** Room F104, Van Steenisgebouw, Einsteinweg 2, 2333 CC Leiden

<https://homepages.cwi.nl/~manegold/DBDM/>

<http://liacs.leidenuniv.nl/~bakkerem2/dbdm/>

DBDM: “Registration”

Please send an email

To: s.manegold@liacs.leidenuniv.nl
Subject: [DBDM-2018] Registration

containing the following information:

- Your full name
- Your email address
- Your student ID
- Your affiliation (university)
- Your program / subject

By Sunday 16 September 2018, 23:59 CEST.

<https://homepages.cwi.nl/~manegold/DBDM/>

<http://liacs.leidenuniv.nl/~bakkerem2/dbdm/>

Databases & Data Mining

Stefan Manegold

Group leader Database Architectures
Centrum Wiskunde & Informatica (CWI)
Amsterdam
<http://homepages.cwi.nl/~manegold/>

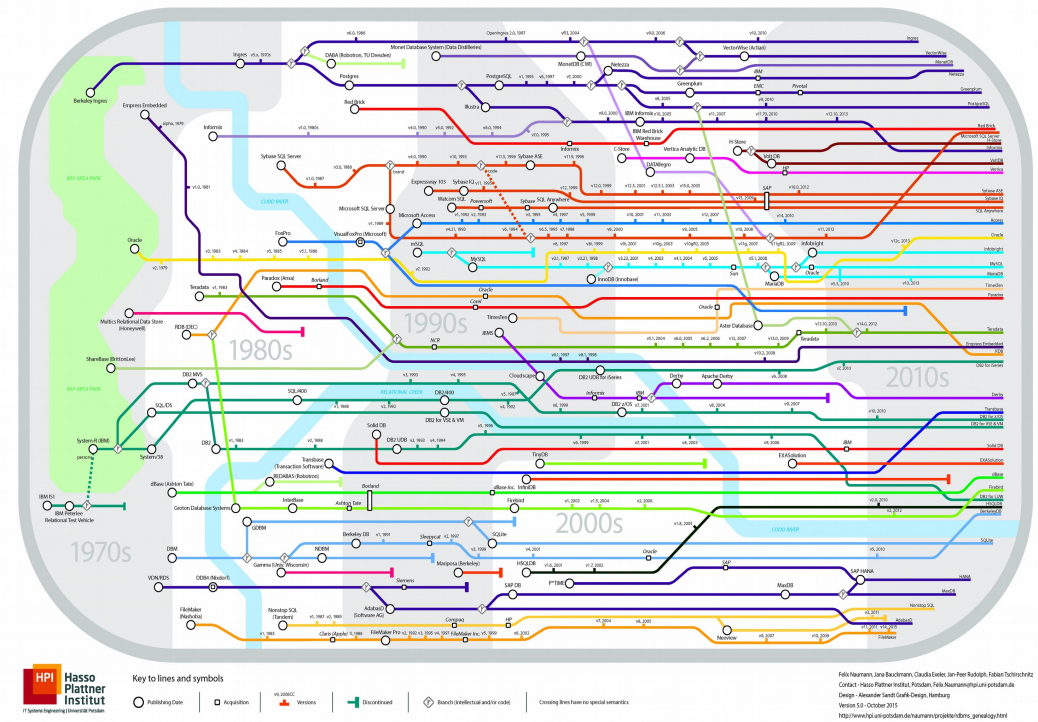
<http://www.monetdb.org/>

Prof. Data Management (0.2 fte)
LIACS & LCDS
Faculty of Science, Leiden University

<https://homepages.cwi.nl/~manegold/DBDM/>

<http://liacs.leidenuniv.nl/~bakkerem2/dbdm/>

Genealogy of Relational Database Management Systems



The age of Big Data



Data



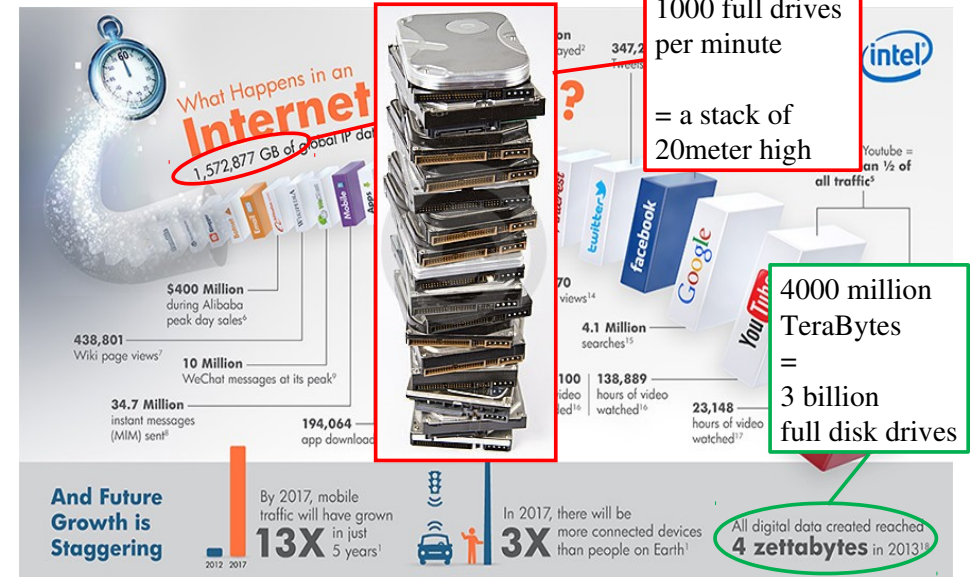
Data Management



Database



Data Mining





Microsoft®

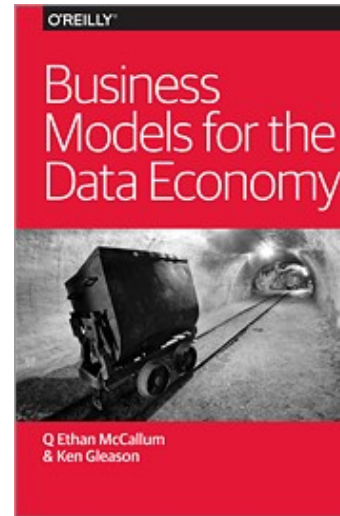


twitter



Google™

facebook

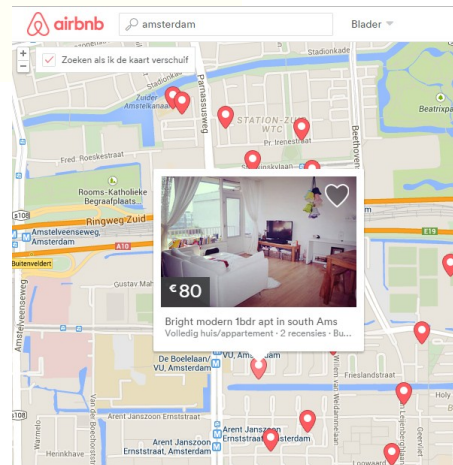
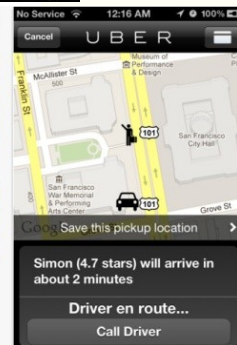
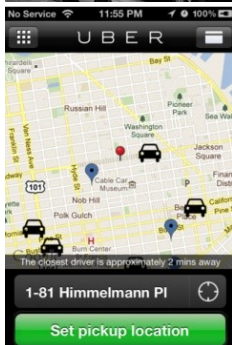


Disruptions by the Data Economy



U B E R

EVERYONE'S PRIVATE DRIVER™



DBDM: Selected Challenges

GIS (LIDAR):

Massive point clouds: 640 Billion (x,y,z) points / 15 TB
=> spatial joins between point cloud and polygons

netherlands
eScience
center

Logistics:



> 5 trillion (10^{12}) GPS points (grows with >60k points/sec)

Seismology:



~ 4 M files, ~ 500 GB (10x compressed)

=> Transparent data ingestion: **Data Vault**

COMMIT/

Remote sensing:



~2 PB satellite image data

=> Array data processing: **SciQL**

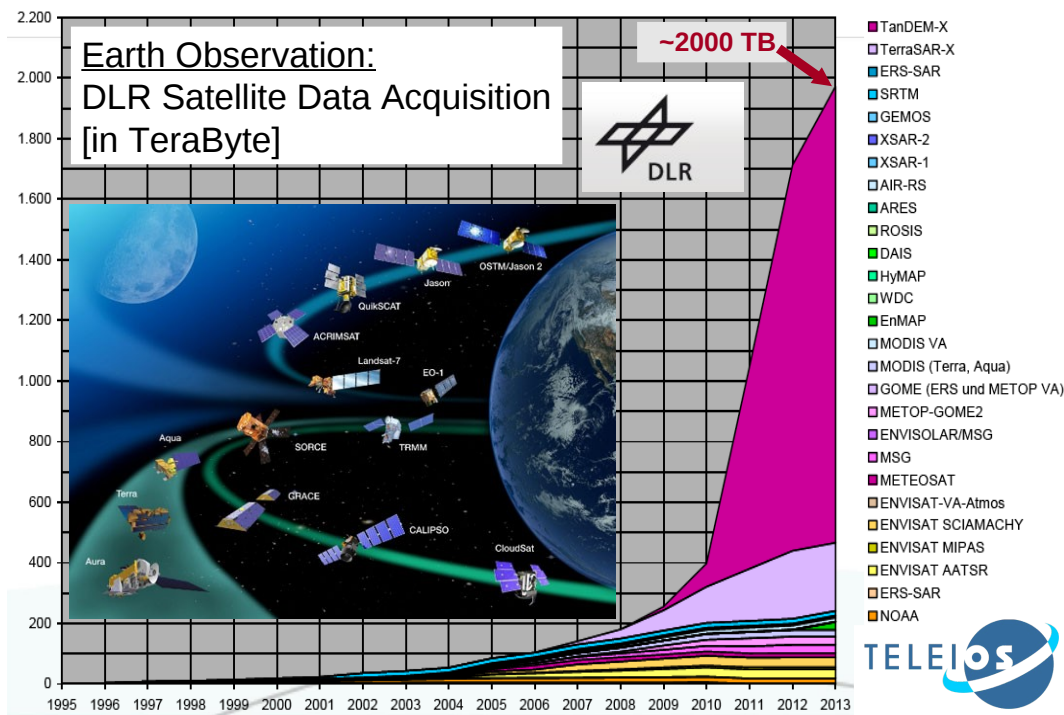


Astronomy:

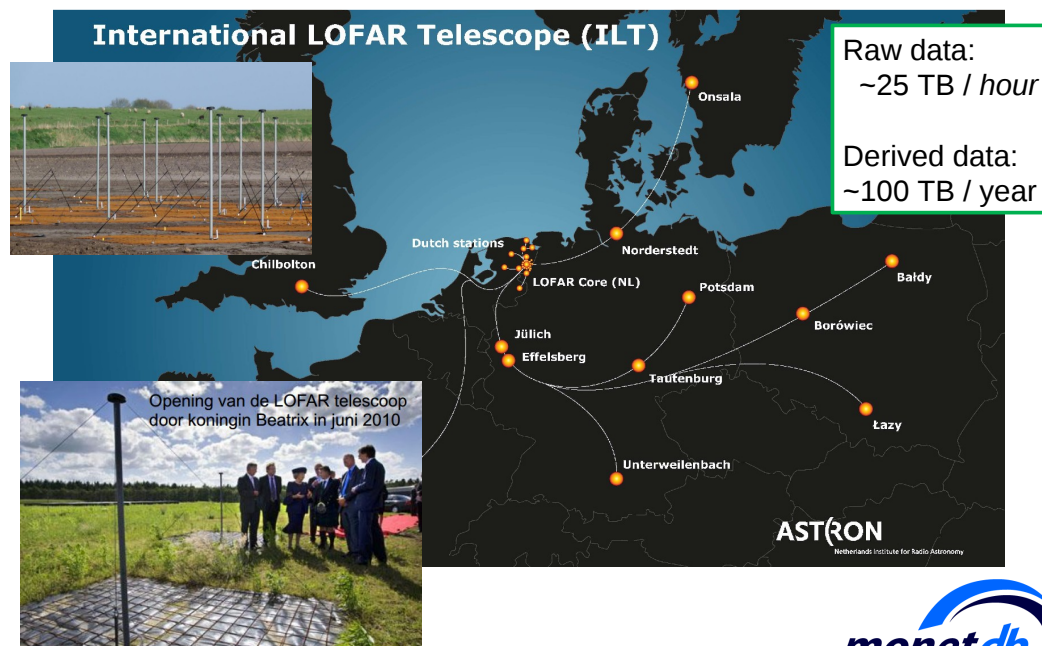


Raw data: 25 TB / hour; derived data: 100 TB / year
=> Transient detection inside DBMS

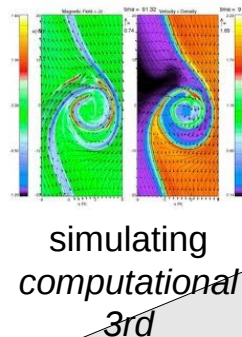
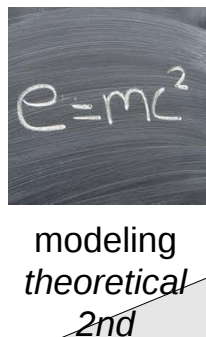
DBDM: Earth Observation



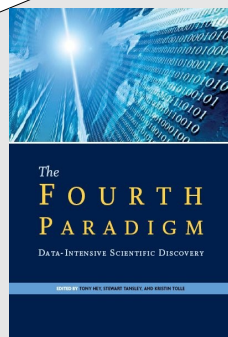
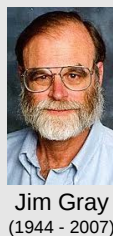
LOFAR Low Frequency Array for Radio Astronomy



Data Disrupting Science: Paradigm Shift in Scientific Research



**collecting &
analyzing data
data exploration
(eScience)**

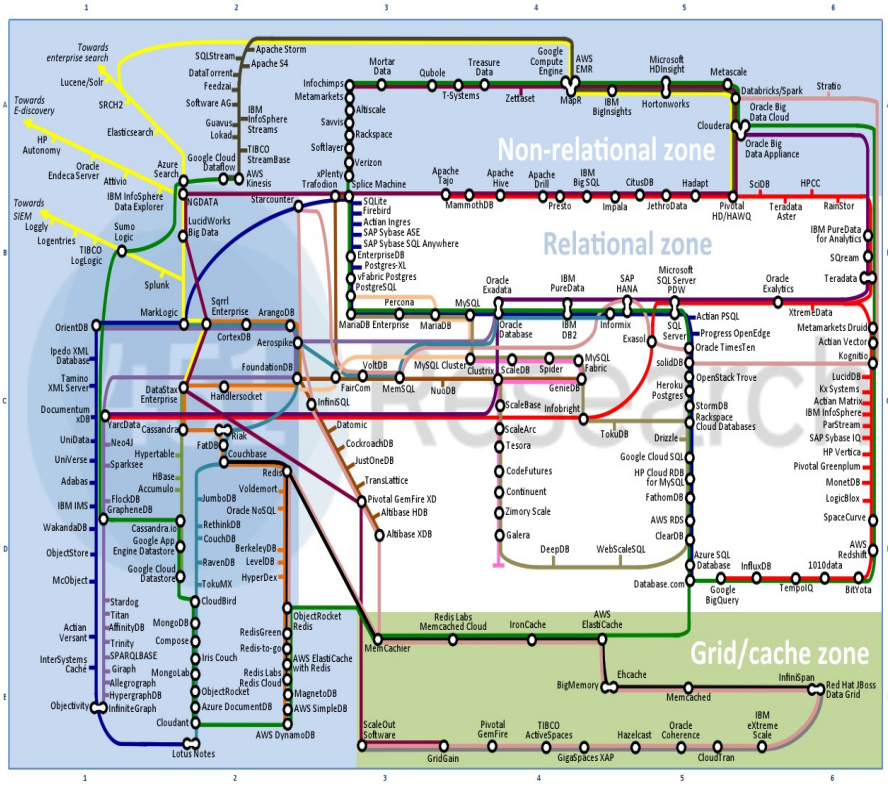


CWI

Database
Architectures

Data Management & Data Mining





451 Research
Data Platforms Map
October 2014

<https://451research.com/dashboard/dpa>

© 2014 by 451 Research LLC.
All rights reserved

BIG DATA LANDSCAPE, VERSION 3.0

Infrastructure

Analytics

Applications

Cross Infrastructure / Open Source

Data Sources

© Matt Turck (@mattturck), Sutan Dong (@sutindong) & FirstMark Capital (@firstmarkcap)